

[Patent Document]

1. Japanese Patent Laid Open
No. 2003-303055

Disk storage system having disk arrays connected with disk adapters through switches

Hitachi, Ltd.

Inventor(s): Tanaka, Katsuya ; Fujimoto, Kazuhisa

Application No. 10/212882, Filed 20020807, A1 Published 20031009

Abstract:

A disk storage system has high throughput between a disk adapter of a disk controller and a disk array. The disk adapter of the disk controller is connected to the disk array through switches. Data on a channel between the switch and a RAID group is multiplexed in the switch to be transferred onto a channel between the switch and the disk adapter and data on the channel between the switch and the disk adapter is demultiplexed in the switch to be transferred onto the channel between the switch and the RAID group. A data transfer rate on the channel between the disk adapter and the switch is made higher than that on the channel.

US.Class: 711114 711154

(19)日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11)特許出願公開番号

特開2003-303055
(P2003-303055A)

(43)公開日 平成15年10月24日(2003.10.24)

| (51)IntCl. ⁷ | G 06 F 3/06 | 国際記号 | F I | ターミナル(参考) |
|-------------------------|-------------|------|-------------|------------------|
| G 06 F 3/06 | 9 01 | | G 06 F 3/06 | 3 01 M 5 B 0 6 5 |
| | 9 02 | | | 3 01 B |
| | 5 4 0 | | | 3 02 A |
| | | | | 5 4 0 |

審査請求 未請求 請求項の版 9 O L (全 14 頁)

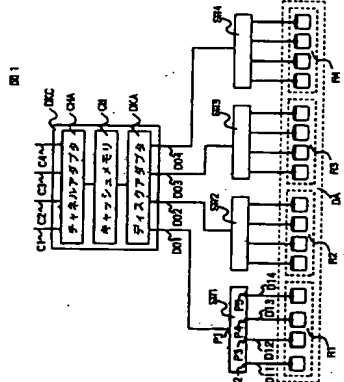
| | | | |
|-----------|-----------------------------|----------|--|
| (21) 出願番号 | 特開2002-106262(P2002-106262) | (71) 出願人 | 000005108 株式会社日立製作所 |
| (22) 出願日 | 平成14年4月9日(2002.4.9) | (72) 発明者 | 株式会社日立製作所 東京都千代田区神田駿河台四丁目6番地 田中 勝也 株式会社日立製作所中央研究所内 東京都国分寺市京産ケ一丁目280番地 株式会社日立製作所中央研究所内 藤本 和久 株式会社日立製作所中央研究所内 東京都国分寺市京産ケ一丁目280番地 株式会社日立製作所中央研究所内 (74) 代理人 100096238 弁護士 伊藤 修 (外1名) Fターム(参考) 5B05 B401 C404 C407 C412 C415 C430 C438 C411 C401 |

(54) 発明の名称 ディスクアダプタとディスクアレイをスイッチを介して接続したディスク装置

(57) 要約

【課題】 ディスクコントローラのディスクアダプタとディスクアレイ間のスループットが高いディスク装置を提供することにある。

【解決手段】 ディスクコントローラ(DKC)のディスクアダプタ(DKA)とディスクアレイ(A)をスイッチ(SW1, SW2, SW3, SW4)を介して接続する。スイッチ(SW1)とRAIDグループ(R1)間のチャネル(C11, D12, D13, D14)上のデータをスイッチ(SW1)において多量化してスイッチ(SW1)とディスクアダプタ(DKA)間のチャネル(C0A)に転送し、スイッチ(SW1)とディスクアダプタ(DKA)間のチャネル(C01)上のデータをスイッチ(SW1)において逆多量化してスイッチ(SW1)とRAIDグループ(R1)間のチャネル(C11, D12, D13, D14)に転送する。ディスクアダプタ(DKA)とスイッチ(SW1)間のチャネル(C01)上のデータ転送速度を、チャネル(C01, D12, D13, D14)のデータ転送速度より高くする。



【特許請求の範囲】

【請求項1】 ディスクコントローラとディスクアレイからなり、前記ディスクコントローラはチャネルアダプタとキャッシュメモリとディスクアダプタを有するディスク装置において、

前記ディスクアダプタと前記ディスクアレイを、バッファメモリを有するスイッチを介して接続し、

前記スイッチは、前記ディスクアダプタが接続されたポートと前記ディスクアレイを構成するディスクドライブが接続された各ポートとの間の接続の切り換えを、入力されたフレーム毎に、該フレーム内の送信優先情報にしたがって行うことを特徴とするディスク装置。

【請求項2】 ディスクコントローラと複数のディスクアレイからなり、前記ディスクコントローラはチャネルアダプタとキャッシュメモリとディスクアダプタを有するディスク装置において、

前記ディスクアレイはループ状に接続した複数のディスクドライブからなり、前記ディスクアダプタと前記複数のディスクアレイとをバッファメモリを有するスイッチを介して接続し、

前記ディスクアダプタと前記スイッチ間のチャネル当りデータ転送速度を、前記スイッチと前記複数のディスクアレイ間のチャネル当りデータ転送速度より高く設定し、

前記スイッチは、前記ディスクアダプタが接続されたポートと前記複数のディスクアレイが接続された各ポートとの間の接続の切り換えを、入力されたフレーム毎に、該フレーム内の送信優先情報にしたがって行うことを特徴とするディスク装置。

【請求項3】 ディスクコントローラとディスクアレイからなり、前記ディスクコントローラはチャネルアダプタとキャッシュメモリとディスクアダプタを有するディスク装置において、

前記ディスクアダプタと前記ディスクアレイを、バッファメモリを有するスイッチを介して接続し、

同一のスイッチに接続したディスクドライブの組み合わせでRAIDグループを構成し、

前記ディスクアダプタと前記スイッチ間のチャネル当りデータ転送速度を、前記スイッチと前記ディスクアレイ間のチャネル当りデータ転送速度より高く設定し、

前記スイッチは、前記ディスクアダプタが接続されたポートと前記RAIDグループを構成するディスクドライブが接続された各ポートとの間の接続の切り換えを、入力されたフレーム毎に、該フレーム内の送信優先情報にしたがって行うことを特徴とするディスク装置。

【請求項4】 第1のディスクコントローラと第2のディスクコントローラと複数のディスクアレイからなり、第1のディスクコントローラは第1のチャネルアダプタ

と第1のキャッシュメモリと第1のディスクアダプタを有し、

第2のディスクコントローラは第2のチャネルアダプタと第2のキャッシュメモリと第2のディスクアダプタを有するディスク装置において、

第1のディスクアダプタと前記複数のディスクアレイとをバッファメモリを有する第1のスイッチを介して接続し、且つ第2のディスクアダプタと前記複数のディスクアレイとをバッファメモリを有する第2のスイッチを介して接続し、さらに第1のスイッチと第2のディスクアダプタを接続し、第2のスイッチと第1のディスクアダプタを接続し、

第1のディスクアダプタと第1のスイッチ間、および第2のディスクアダプタと第1のスイッチ間のチャネル当りデータ転送速度を第1のスイッチと前記複数のディスクアレイ間のチャネル当りデータ転送速度より高く設定し、

第2のディスクアダプタと第2のスイッチ間、および第1のディスクアダプタと第2のスイッチ間のチャネル当りデータ転送速度を第2のスイッチと前記複数のディスクアレイ間のチャネル当りデータ転送速度より高く設定し、

第1のスイッチは、第1のディスクアダプタまたは第2のディスクアダプタが接続されたポートと前記複数のディスクアレイが接続された各ポートとの間の接続の切り換えを、入力されたフレーム毎に、該フレーム内の送信優先情報にしたがって行うことを特徴とするディスク装置。

【請求項5】 第1のディスクコントローラと第2のディスクコントローラと複数のディスクアレイからなり、第1のディスクコントローラは第1のチャネルアダプタと第1のキャッシュメモリと第1のディスクアダプタと第1のキャッシュメモリと第1のディスクアダプタを有するディスク装置において、

第2のディスクコントローラは第2のチャネルアダプタと第2のキャッシュメモリと第2のディスクアダプタと第2のキャッシュメモリと第2のディスクアダプタを有するディスク装置において、

第1のディスクコントローラと第2のディスクコントローラと前記複数のディスクアレイとをバッファメモリを有する第1のスイッチを介して接続し、且つ第2のディスクアダプタと前記複数のディスクアレイとをバッファメモリを有する第2のスイッチを介して接続し、さらに第1のスイッチと第2のディスクアダプタを接続し、第2のスイッチと第1のディスクアダプタを接続し、

第1のディスクアダプタと第1のスイッチ間、および第2のディスクアダプタと第1のスイッチ間のチャネル当りデータ転送速度を第1のスイッチと前記複数のディスクアレイ間のチャネル当りデータ転送速度より高く設定し、

第2のディスクアダプタと第2のスイッチ間、および第1のディスクアダプタと第2のスイッチ間のチャネル当りデータ転送速度を第2のスイッチと前記複数のディスクアレイ間のチャネル当りデータ転送速度より高く設定し、

第1のスイッチは、第1のディスクアダプタまたは第2のディスクアダプタが接続されたポートと前記複数のディスクアレイが接続された各ポートとの間の接続の切り換えを、入力されたフレーム毎に、該フレーム内の送信優先情報にしたがって行うことを特徴とするディスク装置。

【請求項6】 第1のディスクコントローラと第2のディスクコントローラと複数のディスクアレイからなり、第1のディスクコントローラは第1のチャネルアダプタと第1のキャッシュメモリと第1のディスクアダプタと第1のキャッシュメモリと第1のディスクアダプタを有するディスク装置において、

第2のディスクコントローラは第2のチャネルアダプタと第2のキャッシュメモリと第2のディスクアダプタと第2のキャッシュメモリと第2のディスクアダプタを有するディスク装置において、

第1のディスクコントローラと第2のディスクコントローラと前記複数のディスクアレイとをバッファメモリを有する第1のスイッチを介して接続し、且つ第2のディスクアダプタと前記複数のディスクアレイとをバッファメモリを有する第2のスイッチを介して接続し、さらに第1のスイッチと第2のディスクアダプタを接続し、第2のスイッチと第1のディスクアダプタを接続し、

第1のディスクアダプタと第1のスイッチ間、および第2のディスクアダプタと第1のスイッチ間のチャネル当りデータ転送速度を第1のスイッチと前記複数のディスクアレイ間のチャネル当りデータ転送速度より高く設定し、

第2のディスクアダプタと第2のスイッチ間、および第1のディスクアダプタと第2のスイッチ間のチャネル当りデータ転送速度を第2のスイッチと前記複数のディスクアレイ間のチャネル当りデータ転送速度より高く設定し、

第1のスイッチは、第1のディスクアダプタまたは第2のディスクアダプタが接続されたポートと前記複数のディスクアレイが接続された各ポートとの間の接続の切り換えを、入力されたフレーム毎に、該フレーム内の送信優先情報にしたがって行うことを特徴とするディスク装置。

【請求項7】 第1のディスクコントローラと第2のディスクコントローラと複数のディスクアレイからなり、第1のディスクコントローラは第1のチャネルアダプタと第1のキャッシュメモリと第1のディスクアダプタと第1のキャッシュメモリと第1のディスクアダプタを有するディスク装置において、

第2のディスクコントローラは第2のチャネルアダプタと第2のキャッシュメモリと第2のディスクアダプタと第2のキャッシュメモリと第2のディスクアダプタを有するディスク装置において、

第1のディスクコントローラと第2のディスクコントローラと前記複数のディスクアレイとをバッファメモリを有する第1のスイッチを介して接続し、且つ第2のディスクアダプタと前記複数のディスクアレイとをバッファメモリを有する第2のスイッチを介して接続し、さらに第1のスイッチと第2のディスクアダプタを接続し、第2のスイッチと第1のディスクアダプタを接続し、

第1のディスクアダプタと第1のスイッチ間、および第2のディスクアダプタと第1のスイッチ間のチャネル当りデータ転送速度を第1のスイッチと前記複数のディスクアレイ間のチャネル当りデータ転送速度より高く設定し、

第2のディスクアダプタと第2のスイッチ間、および第1のディスクアダプタと第2のスイッチ間のチャネル当りデータ転送速度を第2のスイッチと前記複数のディスクアレイ間のチャネル当りデータ転送速度より高く設定し、

第1のスイッチは、第1のディスクアダプタまたは第2のディスクアダプタが接続されたポートと前記複数のディスクアレイが接続された各ポートとの間の接続の切り換えを、入力されたフレーム毎に、該フレーム内の送信優先情報にしたがって行うことを特徴とするディスク装置。

りデータ転送速度を第1のスイッチと前記複数のディスクアレイ間のチャネル当りデータ転送速度より高く設定し、

第2のディスクアダプタと第2のスイッチ間、および第1のディスクアダプタと第2のスイッチ間のチャネル当りデータ転送速度を第2のスイッチと前記複数のディスクアレイ間のチャネル当りデータ転送速度より高く設定し、

第1のスイッチと第2のスイッチを、第1のディスクアダプタと第2のスイッチ間を接続したチャネルと同等のデータ転送速度を有するチャネルと、第2のディスクアダプタと第1のスイッチ間を接続したチャネルと同等のデータ転送速度を有するチャネルと、を介して接続し、第1のスイッチは、第1のディスクアダプタまたは第2のディスクアダプタまたは第2のスイッチが接続されたポートと前記複数のディスクアレイが接続された各ポートとの間でのポート間の接続の切り換えを、入力されたフレーム毎に、該フレーム内の送信先情報にしたがって行い、

第2のスイッチは、第1のディスクアダプタまたは第2のディスクアダプタまたは第1のスイッチが接続されたポートと前記複数のディスクアレイが接続された各ポートとの間でのポート間の接続の切り換えを、入力されたフレーム毎に、該フレーム内の送信先情報にしたがって行うことを特徴とするディスク装置。

【請求項6】 請求項1乃至請求項5のいずれかの請求項記載のディスク装置において、

前記ディスクアレイからのデータ読み出し時には、前記ディスクアダプタから前記スイッチに転送されるデータを、前記スイッチにおいて逆多量化して前記ディスクアレイに転送することを特徴とするディスク装置。

【請求項7】 請求項1乃至請求項5のいずれかの請求項記載のディスク装置において、

ディスクアダプタからディスクアレイへのデータ書き込み時に、前記ディスクアダプタは、前記ポート間の接続の切り換えが周期的に行われるように、送出するフレームに送信先情報を設定し、

ディスクアレイからディスクアダプタへのデータ読み出し時に、前記スイッチは、ラウンドロビン方式により前記ポート間の接続を切り替えることを特徴とするディスク装置。

【請求項8】 請求項7記載のディスク装置において、周期的に切り換えられるポート数を、ディスクアダプタとスイッチ間のチャネル当りデータ転送速度の、スイッチとディスクアレイ間のチャネル当りデータ転送速度に対する比、と同程度に設定することを特徴とするディスク装置。

図。
【請求項9】 請求項1乃至請求項5のいずれかの請求項記載のディスク装置において、

前記ディスクアダプタと前記スイッチ間を光ファイバケーブルで接続し、前記スイッチと前記ディスクアレイ間をメタルケーブルで接続することを特徴とするディスク装置。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】 本発明は、コンピュータシステムにおける2次記憶装置に関し、特に入出力データ転送性能が高いディスク装置に関する。

【0002】

【従来の技術】 現在のコンピュータシステムにおいては、CPU（中央処理装置）が必要とするデータは2次記憶装置に保存され、CPUなどが必要とするときに必じて2次記憶装置に対してデータの書き込みおよび読み出しを行う。この2次記憶装置としては、一般に不揮発性記憶媒体が使用され、代表的なものとして磁気ディスク装置や、光ディスクなどのディスク装置がある。近年高度情報化に伴い、コンピュータシステムにおいて、この種の2次記憶装置の高性能化が要求されている。

【0003】 図9に、従来のディスク装置のブロック図を示す。図9において、ディスク装置はディスクコントローラDKCとディスクアレイDAで構成される。ディスクコントローラDKCは、上位側CPU（図示せず）とディスク装置を接続するチャネルアダプタCHAと、ディスクアレイDAに対して読み書きするデータを一時保存するキャッシュメモリCMと、ディスクコントローラDKCとディスクアレイDAを接続するディスクアダプタDKAからなる。チャネルアダプタCHAとキャッシュメモリCMとディスクアダプタDKAは、バスまたはスイッチで相互接続されている。チャネルアダプタCHAはC1、C2、C3、C4の4本のチャネルでCPUと接続している。ディスクアダプタDKAはD1、D2、D3、D4の4本のチャネルでディスクアレイと接続している。ここでディスクアレイDAはディスクグループR1、R2、R3、R4からなり、それぞれRAIDグループを構成する。

【0004】 チャネルC1、C2、C3、C4から入力された書き込みデータは、キャッシュメモリCMに敵データを書き込みむと同時に、敵データをブロックサイズ単位に分割し、チャネルD1、D2、D3、D4の4本のチャネルには前記分割データから計算したパリティを、ディスクアダプタDKAからディスクアレイDAへ送出する。データ読み出し時は、まずキャッシュメモリCM内に敵データの有無を調べる。有る場合は、キャッシュメモリCMからディスクアレイDAへデータを転送し、ディスクアレイDAからディスクアダプタDKAへデータを転送し、ディスクアダプタDKAからキャッシュメモリCMへデータを転送し、キャッシュメモリCMからディスクコントローラDKCへデータを転送し、ディスクコントローラDKCからCPUへデータを転送する。

【0005】 図10に、本発明のディスク装置のブロック図を示す。図10において、ディスク装置はディスクコントローラDKCとディスクアレイDAで構成される。ディスクコントローラDKCは、上位側CPU（図示せず）とディスク装置を接続するチャネルアダプタCHAと、ディスクアレイDAに対して読み書きするデータを一時保存するキャッシュメモリCMと、ディスクコントローラDKCとディスクアレイDAを接続するディスクアダプタDKAからなる。チャネルアダプタCHAとキャッシュメモリCMとディスクアダプタDKAは、バスまたはスイッチで相互接続されている。チャネルアダプタCHAはC1、C2、C3、C4の4本のチャネルでCPUと接続している。ディスクアダプタDKAはD1、D2、D3、D4の4本のチャネルでディスクアレイと接続している。ここでディスクアレイDAはディスクグループR1、R2、R3、R4からなり、それぞれRAIDグループを構成する。

【0006】 チャネルC1、C2、C3、C4から入力された書き込みデータは、キャッシュメモリCMに敵データを書き込みむと同時に、敵データをブロックサイズ単位に分割し、チャネルD1、D2、D3、D4の4本のチャネルには前記分割データから計算したパリティを、ディスクアダプタDKAからディスクアレイDAへ送出する。データ読み出し時は、まずキャッシュメモリCM内に敵データの有無を調べる。有る場合は、キャッシュメモリCMからディスクアレイDAへデータを転送し、ディスクアダプタDKAからキャッシュメモリCMへデータを転送し、キャッシュメモリCMからディスクコントローラDKCへデータを転送し、ディスクコントローラDKCからCPUへデータを転送する。

【0007】 図11に、本発明のディスク装置のブロック図を示す。図11において、ディスク装置はディスクコントローラDKCとディスクアレイDAで構成される。ディスクコントローラDKCは、上位側CPU（図示せず）とディスク装置を接続するチャネルアダプタCHAと、ディスクアレイDAに対して読み書きするデータを一時保存するキャッシュメモリCMと、ディスクコントローラDKCとディスクアレイDAを接続するディスクアダプタDKAからなる。チャネルアダプタCHAとキャッシュメモリCMとディスクアダプタDKAは、バスまたはスイッチで相互接続されている。チャネルアダプタCHAはC1、C2、C3、C4の4本のチャネルでCPUと接続している。ディスクアダプタDKAはD1、D2、D3、D4の4本のチャネルでディスクアレイと接続している。ここでディスクアレイDAはディスクグループR1、R2、R3、R4からなり、それぞれRAIDグループを構成する。

【0008】 チャネルC1、C2、C3、C4から入力された書き込みデータは、キャッシュメモリCMに敵データを書き込みむと同時に、敵データをブロックサイズ単位に分割し、チャネルD1、D2、D3、D4の4本のチャネルには前記分割データから計算したパリティを、ディスクアダプタDKAからディスクアレイDAへ送出する。データ読み出し時は、まずキャッシュメモリCM内に敵データの有無を調べる。有る場合は、キャッシュメモリCMからディスクアレイDAへデータを転送し、ディスクアダプタDKAからキャッシュメモリCMへデータを転送し、キャッシュメモリCMからディスクコントローラDKCへデータを転送し、ディスクコントローラDKCからCPUへデータを転送する。

【0009】 図12に、本発明のディスク装置のブロック図を示す。図12において、ディスク装置はディスクコントローラDKCとディスクアレイDAで構成される。ディスクコントローラDKCは、上位側CPU（図示せず）とディスク装置を接続するチャネルアダプタCHAと、ディスクアレイDAに対して読み書きするデータを一時保存するキャッシュメモリCMと、ディスクコントローラDKCとディスクアレイDAを接続するディスクアダプタDKAからなる。チャネルアダプタCHAとキャッシュメモリCMとディスクアダプタDKAは、バスまたはスイッチで相互接続されている。チャネルアダプタCHAはC1、C2、C3、C4の4本のチャネルでCPUと接続している。ディスクアダプタDKAはD1、D2、D3、D4の4本のチャネルでディスクアレイと接続している。ここでディスクアレイDAはディスクグループR1、R2、R3、R4からなり、それぞれRAIDグループを構成する。

【0010】 チャネルC1、C2、C3、C4から入力された書き込みデータは、キャッシュメモリCMに敵データを書き込みむと同時に、敵データをブロックサイズ単位に分割し、チャネルD1、D2、D3、D4の4本のチャネルには前記分割データから計算したパリティを、ディスクアダプタDKAからディスクアレイDAへ送出する。データ読み出し時は、まずキャッシュメモリCM内に敵データの有無を調べる。有る場合は、キャッシュメモリCMからディスクアレイDAへデータを転送し、ディスクアダプタDKAからキャッシュメモリCMへデータを転送し、キャッシュメモリCMからディスクコントローラDKCへデータを転送し、ディスクコントローラDKCからCPUへデータを転送する。

が、ディスクアダプタのポート数増加は制御を複雑にする。第2の従来技術では、ディスクアダプタとディスクアレイとの間にスイッチを適用することによりディスク増設ポート数を増加させることができるが、チャネル当りのデータ転送速度はディスクアレイのデータ転送速度に制限されるので、ディスクアダプタとディスクアレイ間のスループットが性能ネックになるという問題があった。第3の従来技術は、ディスクの回転待ち時間の影響を低減できる技術であり、フロントエンドとバックエンドのスループット増強は低減できないという問題があった。

【0008】 本発明の目的は、ディスクアダプタとディスクアレイ間のスループットが高いディスク装置を提供することにある。本発明の他の目的は、ディスクアダプタとディスクアレイ間のスループットが高く、且つディスクドライブ接続台数が多いディスク装置を提供することにある。本発明のさらに他の目的は、信頼性が高いディスクアレイを有するディスク装置を提供することである。本発明のさらに他の目的は、信頼性が高いディスクアダプタとディスクアレイ間ネットワークを有するディスク装置を提供することにある。本発明のさらに他の目的は、ディスクからの読み出しおよびディスクへの書き込みを高スループット化できるディスク装置を提供することにあり、本発明のさらに他の目的は、高スループットを維持できるディスク装置を提供することである。本発明のさらに他の目的は、高スループットで低コストなディスク装置を提供することである。

【0009】

【課題を解決するための手段】 上記目的を達成するため、本発明は、ディスクコントローラとディスクアレイからなり、ディスクコントローラはチャネルアダプタとキャッシュメモリとディスクアダプタを有するディスク装置であり、ディスクアダプタとディスクアレイを、バスを介して接続し、ディスクアレイ間のチャネル当りデータを、スイッチとスイッチ間のチャネル当りデータ転送速度を、スイッチとスイッチ間のチャネル当りデータ転送速度より高く設定し、スイッチは、ディスクアダプタが接続されたポートとディスクアレイを構成するディスクアレイが接続された各ポートとの間でのポート間の接続の切り換えを、入力されたフレーム毎に、該フレーム内の送信先情報にしたがって行っている。また、前記ディスクアレイはグループ化された複数のディスクドライブからなり、前記ディスクアダプタと前記複数のディスクアレイとをバスファームを有するスイッチを介して接続し、ディスクアダプタとスイッチ間のチャネル当りデータ転送速度を、スイッチと複数のディスクアレイ間のチャネル当りデータ転送速度より高く設定し、

【0010】

【請求項1】 請求項1乃至請求項5のいずれかの請求項記載のディスク装置において、

ディスクアダプタからディスクアレイへのデータ書き込み時に、前記ディスクアダプタは、前記ポート間の接続の切り換えが周期的に行われるように、送出するフレームに送信先情報を設定し、

ディスクアレイからディスクアダプタへのデータ読み出し時に、前記スイッチは、ラウンドロビン方式により前記ポート間の接続を切り替えることを特徴とするディスク装置。

【請求項2】 請求項1記載のディスク装置において、周期的に切り換えられるポート数を、ディスクアダプタとスイッチ間のチャネル当りデータ転送速度の、スイッチと複数のディスクアレイ間のチャネル当りデータ転送速度に対する比、と同程度に設定することを特徴とするディスク装置。

スイッチは、ディスクアダプタが接続されたポートと複数のディスクアレイが接続された各ポートとの間でのポート間の接続の切り換えを、入力されたフレーム毎に、破フレーム内の送信優先権にしたがって行っている。また、前記ディスクアダプタと前記ディスクアレイを、バックアップメモリを有するスイッチを介して接続し、同一のスイッチに接続したディスクドライブの組み合わせでR A I Dグループを構成し、ディスクアダプタとスイッチ間のチャネル当りデータ転送速度を、スイッチとディスクアレイ間のチャネル当りデータ転送速度より高く設定し、スイッチは、ディスクアダプタが接続されたポートとR A I Dグループを構成するディスクドライブが接続された各ポートとの間でのポート間の接続の切り換えを、入力されたフレーム毎に、破フレーム内の送信優先権にしたがって行っている。また、第1のディスクコントローラと第2のディスクコントローラと複数のディスクアレイからなり、第1のディスクコントローラは第1のチャネルアダプタと第1のキャッシュメモリと第1のディスクアダプタを有し、第2のディスクコントローラは第2のチャネルアダプタと第2のキャッシュメモリと第1のディスクアダプタを有する第1のキャッシュメモリと第2のディスクアダプタを有する第2のキャッシュメモリとをバックアップメモリを有する第1のスイッチを介して接続し、且つ第2のディスクアダプタと前記複数のディスクアレイとをバックアップメモリを有する第2のスイッチを介して接続し、さらに第1のスイッチと第2のディスクアダプタを接続し、第2のスイッチと第1のディスクアダプタを接続し、第2のディスクアダプタと第2のスイッチと、および第1のディスクアダプタと第2のスイッチ間のチャネル当りデータ転送速度を第2のスイッチと前記複数のディスクアレイ間のチャネル当りデータ転送速度より高く設定し、第1のスイッチは、第1のディスクアダプタまたは第2のディスクアレイが接続されたポートと前記複数のディスクアレイが接続された各ポートとの間でのポート間の接続の切り換えを、入力されたフレーム毎に、破フレーム内の送信優先権にしたがって行っている。第2のスイッチは、第1のディスクアダプタまたは第2のディスクアダプタが接続されたポートと前記複数のディスクアレイが接続された各ポートとの間でのポート間の接続の切り換えを、入力されたフレーム毎に、破フレーム内の送信優先権にしたがって行っている。また、さらに、上記第1のスイッチと第2のスイッチと、上記第1のディスクアダプタと第2のスイッチ間を接続したチャネルと同等のデータ転送速度を有するチャネルと、第2のディスクアダプタと第1のスイッチ間を接続したチャネルと同等のデータ転送速度を有するチャネルとを、介して接続している。また、前記ディスクアレイからのデータ転送に出る際には、前記ディスクアレイから前記スイッチに転送されるデータを前記スイッチにおいて多重化して前記ディスクアダプタに転送し、前記ディ

スクアレレイへのデータ書き込み時には、前記ディスプレイから前記スイッチに転送されるデータを前記ディスプレイにおいて逆多重化し、前記ディスプレイから前記スクアレレイへ送出する。また、前記ディスプレイからのデータ書き込み時に、前記ディスプレイのデータ書き込み時に、前記ディスプレイ間の接続の切り替えが行われ、前に送出するフレームに送信先情報を設定し、前に送り出すときに、前記スイッチは、ラウンドロビン方式により前記ディスプレイ間の接続を切り替えるようにしている。

また、さらに、切り替えるポート数を、ディスプレイアダプタとスロット間のチャネル当りデータ転送速度の、スイッチとディスプレイ間のチャネル当りデータ転送速度に対する比、と同程度に設定している。また、前記ディスプレイ間と前記スイッチ間の光ファイバケーブルをメタルケーブルで接続するようにしている。

【00101】
 〔発明の実施の形態〕以下、図面を参照して本発明の実施の形態を詳細に説明する。図1に本発明の、第1の実施の形態であるディस्क装置の構成を示す。本実施の形態のディस्क装置は、ディスクコントローラDKCとディスクアレイDAからなる。ディスクコントローラDKCは、チャネルアダプタCHAと、キャッシュメモリCMと、ディスクアダプタDKAからなる。チャネルアダプタCHAは、上位CPU（図示せず）とディスクコントローラDKCとがデータを送受信する際の制御を行う。C1、C2、C3およびC4は、チャネルアダプタDKAがCPUと通信するチャネルである。キャッシュメモリCMは、本実施の形態のディस्क装置に入出力するデータを一時保持するメモリである。ディस्कアダプタDKAは、ディスクコントローラDKCとディスクアレイDAとがデータを送受信する際の制御を行う。ディスクアダプタDKAは、チャネルD01、D02、D03、D04を介して、ディスクアレイDAと接続する。ディスクアダプタDKAとディスクアレイDAは、チャネルD01、D02、D03、D04上で全二重通信が可能である。

【0011】ここで、本実施の形態のディスク装置は、ディスクアダプタDKAとディスクアレイDAを、スイッチSW1、SW2、SW3、SW4を介して接続している点に特徴がある。ディスクアレイDAは、ディスクグループR1、R2、R3、R4からなる。ディスクグループR1は、スイッチSW1介してディスクアダプタDKAと接続する。同様に、ディスクグループR2はスイッチSW2を介して、ディスクグループR3はスイッチSW3を介して、ディスクグループR4はスイッチSW4介して、それぞれディスクアダプタDKAと接続する。

【0012】本実施の形態のディスク装置においてRA

IDシステムを構築する場合は、ディスクグループR1、R2、R3、R4を、それぞれRAIDグループとして構築する。本実施の形態では、4個のディスクタイプでRAIDグループを構成しているが、RAIDグループを構成するドライブ数を4個に限るものではない。各ディスクグループへのデータ読み出しまたは書き込み時のデータの流れを、ディスクグループR1を例にして述べる。ここでR1はRAIDレベル5のRAIDグループである。チャネルC1、C2、C3、C4からディスクグループR1へ書き込むためにCPIから送信されたデータは、ディスクアダプタDKAにおいてブロック単位に分割されると同時に、該ブロック単位に分割されたデータからパリティが生成される。該ブロック単位に分割されたデータと、生成されたパリティは、チャネルD0、D1を通りスイッチSW1へ入力される。スイッチSW1は、RAID制御に伴い、該ブロック単位に分割されたデータと、生成されたパリティとをローテーションし、チャネルD11、D12、D13、D14へ分配する。データ読み出し時は、ディスクアダプタDKAは、D1、D12、D13、D14を介してディスクグループR1からブロック単位に分割されたデータを読み出し、スイッチSW1でシリアル化して、チャネルD01を通じて読み出しデータを受信する。

[0013] 図9に示した従来のディスク装置では、ディスクアダプタDKAに接続したチャネルD1、D2、D3、D4上で、既にディスクレイブへの書き込みデータおよびパリティが別々のチャネルに分配されていた。それに対し、本実施の形態のディスク装置においては、スイッチSW1通過後に別々のチャネルに分配される点から従来と異なる。

【0014】次に、本実施形態のディスク装置の特性であるスライツの動作を、スライツSW1を例にとり説明する。SW2～SW4の動作もSW1の動作と同様である。図1に示すように、スライツSW1は出力ポートP1、P2、P3、P4、P5を有する。ポート1、2、3、4、P5は、全二重通信可能な出力ポートであり、ポート毎にバッファメモリを有している。スライツSW1の内部構成を図2と図3に示す。簡単のため、データの通行方向によりスライツ動作を分けて説明する。また、チャネルD01、D11、D12、D13、D14上を通れるデータは、フレーム単位で送受信され、かつデータは8B10B変換で符号化されてい

【0015】図2は、ポートP1からブロック内のフレームを入力し、ポートP2、P3、P4、P5から出力する場合を示す。これはディスプレイへの書き込み時のスイッチ動作に相当する。スイッチSW1は図2に示すように、クロスバスイッチXSWと、スイッチコントローラCTLからなる。クロスバスイッチXSWは5×5のクロスバスイッチであり、入力ポートin1、in2、in3、in4、in5、出力ポートout1、out2、out3、out4、out5を有する。

2、ln3、ln4、ln5と、出力ポートout1、out2、out3、out4、out5を有する。ポートP1から入力したフレームは、シリアルパラレル変換回路SP1と、バッファメモリBM1と、8B10B変換回路DEC1を經由し、スイッチコントローラCTLと入力ポートln1へ入力される。スイッチコントローラCTLにおいて、入力フレームのヘッダ部分に記された送信元アドレスを解釈し、クロスマスイッチXSWを切り換える。例として、ポートP2が出力先として選ばれた場合は、入力したフレームは出力ポートout2と、8B10B変換回路ENC2と、バッファメモリBM2と、パラレルシリアル変換回路PS2を經由し、ポートP2から出力される。ここで、バッファメモリBM1、BM2はFIFO(First-In-First-Out)メモリである。

【0016】シリアルパラレル変換装置SP1は、8B10B符号化されたシリアルデータを10b11個のパラレルデータに変換し、ポートP1におけるデータ転送速度の1/10の速度に同期してパッファメモリBM1に書き込む。8B10BデコードDEC1は、クロスバスイッチXSWの動作速度に同期して、10b11パラレルデータをパッファメモリBM1から読み出し、8B10B符号化して、8b11パラレルデータに変換する。8B10BエンコードENC2は、クロスバスイッチXSWでスイッチされた8b11パラレルデータを再び8B10B符号化し、10b11パラレルデータに変換し、ポートP2に出力する。以上によりスイッチSW1は、ポートP1におけるデータ転送速度からポートP2におけるデータ転送速度へ速度変換する。

【0017】図4は、ポートP1へ入力するフレームと、ポートP2、P3、P4、P5から出力されるフレームを示した図である。波形の凸はフレームが存在する時間、凹はフレームが存在していない時間を示している。フレームは伝送するデータ容量に依らず一定のフレーム長が変化するが、ここではディスクアレイへのシリアルアクセスが行われており、フレーム長が一定である。図4では、入力ポートP1でのデータ転送速度が出力ポートP2、P3、P4、P5におけるデータ転送速度のm倍あるとする。従って、ポートP1におけるフレームm倍あるとする。従って、ポートP1における時間T3へ伸びている。ここで $T3 = m \times T1$ である。

【0.0.18】入力のデータ転送速度が速く、且つ出力のデータ転送速度が遅い場合は、スイッチを定期的に切り換えないと出力ポートのパufferメモリが溢れ、スルーブットが低下する。フレイムがスルーブットの低下無く

スイッチを通ずるには、図4のように周期的に出力ポートを切り替える必要がある。スイッチ切り替えポート数を n とすると、スイッチ切り替え周期 $T = n \times T1$ である。フレームの無い時間は無視した。 $T2 \geq T3$ ならば、フレームの面欠無く、スループットの低下は起こらない。 $T2 \geq T3$ は n と同じである。つまり、ディスクアレイへのデータを送る時に、スイッチにおいてスループット低下を起こさないための条件は、周期的に切り替えるスイッチポート数 n を、ディスクアダプタとスイッチ間のチャネル当りデータ転送速度の、スイッチとディスクアレイ間のチャネル当りデータ転送速度に対する比 m 、以上に設定することである。この条件が保たれれば、スイッチSW1は、ポートP1から入力したデータをバッファメモリにおいて速度変換し、フレーム単位で周期的に切り替えることにより逆多重化し、ポートP2、P3、P4、P5へ分配して出力する。スイッチを、周期的に切り替える方法の一つは、スイッチに接続した、ディスクグループをRAIDグループとすることである。RAIDのストライピング原理に従えば、スイッチは周期的に切り替わる。

[0019] 図3は、ポートP2、P3、P4、P5からフレームを入力し、ポートP1から出力する場合を示す。これはディスクアレイからの読み出し時のスイッチ動作に相当する。例えば、ポートP2から入力したフレームは、シリアルパラレル変換装置SP2と、バッファメモリBM2と、8B10B変換デコードDEC2を経由し、スイッチコントローラCTLと入力ポートIn2へ入力される。スイッチコントローラCTLにおいて、入力フレームのヘッダ部分に書かれた送信先アドレスを解読し、クロスバスイッチXSWを切り替える。図3の場合、ラウンドロビン方式によりクロスバスイッチXSWを切り替えて、順番にポートP2、P3、P4、P5から入力されるデータは全てポートP1へ出力する。すなわち、読み出し時は、複数の入力ポート(P2、P3、P4、P5)に同時にフレームが届く。これら複数の入力フレームは同時に出力ポートに届く必要はない。スイッチは、総当たり的に入力ポート間接続を切り替えることにより、これら複数の入力フレームをフレームずつ出力ポート(P1)へ転送する。このように、スイッチを総当たり的に切り替える方式を、ラウンドロビン(Round Robin)方式と呼ぶ。ラウンドロビン方式により、結果的にスイッチは周期的に切り替わることになる。なお、読み出し時においても、スイッチはフレーム内送信先情報に従って切り替わることに違いない。フレームは入力ポートOut1と、8B10B変換エンコードDEC1と、バッファメモリBM1と、パラレルシリアル変換装置PS1を經由して、ポートP1から出力される。

[0020] シリアルパラレル変換装置SP2は、8B10B符号化されたシリアルデータを10ビット幅のバ

ラレルデータに変換し、ポートP2におけるデータ転送速度の1/10の速度に同期してバッファメモリBM2に書き込む。8B10BデコードDEC2は、クロスバスイッチXSWの動作速度に同期して、10ビットパラレルデータをバッファメモリBM2から読み出し、8B10B符号化して、8ビットパラレルデータに変換する。8B10BエンコードENC1は、クロスバスイッチXSWでスイッチされた8ビットパラレルデータを再び8B10B符号化し、10ビットパラレルデータに変換後、クロスバスイッチXSWの動作速度に同期してバッファメモリBM1に書き込む。パラレルシリアル変換装置PS1は、ポートP1におけるデータ転送速度の1/10の速度に同期して、10ビットパラレルデータをバッファメモリBM1から読み出し、シリアル化して、ポートP1から出力する。以上によりスイッチSW1は、ポートP2におけるデータ転送速度からポートP1におけるデータ転送速度へ速度変換する。

[0021] 図5は、ポートP2、P3、P4、P5へ入力するフレームと、ポートP1から出力されるフレームを示した図である。波形の凸はフレームが存在する時間、凹はフレームが存在していない時間を示している。フレームは伝送するデータ容量に従ってそのフレーム長が変化するが、ここではディスクアレイへのシーケンシャルアクセスが行われており、フレーム長が一定である。図5では、入力ポートP1でのデータ転送速度が出力ポートP2、P3、P4、P5におけるデータ転送速度の m 倍あるとする。従って、ポートP5におけるフレームFe5の時間T4は、ポートP1からの出力時にT5へ縮んでいる。ここで $T4 = m \times T5$ である。フレームFe2、Fe3、Fe4、Fe5をポートP1から出力するのにかかる時間をT6とする。スイッチ切り替えポート数を n とすると、 $T6 = n \times T5$ である(フレームの無い時間は無視した)。スイッチにおいて幅域によるスループット低下を防止するためには、 $T6 \leq T4$ とする必要がある。 $T6 \leq T4$ は $n \leq m$ と同じである。

[0022] つまり、ディスクアレイからのデータ読み出し時に、スイッチにおいてスループット低下を起こさないための条件は、周期的に切り替えるスイッチポート数 n を、ディスクアダプタとスイッチ間のチャネル当りデータ転送速度の、スイッチとディスクアレイ間のチャネル当りデータ転送速度に対する比 m 、以下に設定することである。この条件が保たれれば、スイッチSW1は、ポートP2、P3、P4、P5から入力したデータをバッファメモリにおいて速度変換し、フレーム単位で周期的に切り替えることにより逆多重化し、ポートP1へ出力する。よって、ディスクアレイへの書き込みおよびディスクアレイからの読み出しを高スループット化するために、 $n \leq m$ 、つまり、周期的に切り替えるポート数を、ディスクアダプタとスイッチ間のチャネル当りデータ転送速度の、スイッチとディスクアレイ間のチャ

ネル当りデータ転送速度に対する比、と同程度に設定すればよいことが分かる。

[0023] 例えば、ディスクアダプタとスイッチ間の4Gbpsのチャネル1本で接続し、スイッチとディスクアレイ間を1Gbpsのチャネル4本で接続する。また、ディスクアダプタとスイッチ間の10Gbpsのチャネル1本で接続し、スイッチとディスクアレイ間を2Gbpsのチャネル4本で接続する。この場合、スイッチ入出力ポート間でスループットのバランスが取れないので、実効的なスループットは $2Gbps \times 4 = 8Gbps$ とする。

[0024] 以上より、スイッチSW1において速度変換と多重化、逆多重化が行われるので、チャネルD1、D12、D13、D14上のデータ転送速度が低速でも、チャネルD01、D02、D03、D04でのデータ転送速度は高速にできる。つまり、ディスクアダプタDKAとディスクアレイDA間のスループットを上向きにできる。本実施の形態のディスク装置におけるデータ転送方式としては、ファイバチャネルやインフィニバンドが使用できる。

[0025] 図6は、第1の実施の形態のディスク装置において、ディスクドライブの増設方法を示した図である。図6では図1に対して、ディスクグループR5とR6が増設されている。ディスクドライブを増設する際に、スイッチSW1とSW2としてポート数の多いスイッチを使用している。ディスクドライブを増設すると、スイッチのディスクアレイ側のスループットが増加し、ディスクアダプタ側のスループットバランスが崩れるので、スイッチの速度変換機能が有効に働かなくなる可能性がある。そこでスイッチSW1では、ディスクアダプタDKAとの間に、新規チャネルD05を増設している。また、スイッチSW2の場合は新規チャネルを増設せず、チャネルD02の信号伝送速度を増加させることで、ディスクアダプタ側とディスクアレイ側のスループットバランスを取っている。例えばスイッチSW1では、スイッチとディスクアレイ間を1Gbpsのチャネル8本で接続し、ディスクアダプタとスイッチ間を4Gbpsのチャネル2本で接続する。スイッチSW2では、スイッチとディスクアレイ間を1Gbpsのチャネル8本で接続し、ディスクアダプタとスイッチ間を10Gbpsのチャネル1本で接続する。このように、本実施の形態のディスク装置は、スイッチのポート数に応じたドライブ増設方法は、1ポート当たり接続できるドライブ数が少ないATA(AT Attachment)方式ディスクドライブを増設するのに適用できる。

[0026] 図7に本発明の、第2の実施の形態であるディスク装置の構成を示す。本実施の形態のディスク装置は、第1の実施の形態のディスク装置に対して、ディスクアレイ部分の構成方法が異なる。本実施の形態のデ

ィスク装置は、ディスクコントローラDKCと、4個のディスクアレイDA1、DA2、DA3、DA4からなる。ディスクコントローラDKCは、チャネルアダプタCHA、キャッシュメモリCM、ディスクアダプタDKAからなる。ディスクアレイDA1とディスクアダプタDKAは、チャネルD01とスイッチSW1を介して接続する。同様に、ディスクアレイDA2はチャネルD02とスイッチSW2を介して、ディスクアレイDA3はチャネルD03とスイッチSW3を介して、ディスクアレイDA4はチャネルD04とスイッチSW4を介して、それぞれディスクアダプタDKAと接続する。スイッチSW1、SW2、SW3とSW4は、第1の実施の形態と同様に速度変換と多重化、逆多重化を行うスイッチとして機能する。本実施の形態におけるディスクアダプタDKAと、スイッチSW1、SW2、SW3、SW4と、ディスクアレイDA1、DA2、DA3、DA4との間のデータ転送方式は、ファイバチャネルを使用している。スイッチSW1、SW2、SW3、SW4はファイバチャネルスイッチである。

[0027] 本実施の形態におけるディスクアレイの構成を、ディスクアレイDA1を例に述べる。ディスクアレイDA1、DA2、DA3、DA4は、同様のドライブ構成である。ディスクアレイDA1は、チャネルD11上に接続した4個のディスクからなるディスクアレイと、D12上に接続した4個のディスクからなるディスクアレイと、D13上に接続した4個のディスクからなるディスクアレイと、D14上に接続した4個のディスクからなるディスクアレイと、からなる。チャネルD11を例にとると、ディスクドライブDK1、DK2、DK3、DK4が、チャネルD11上に接続されている。このように、多数のドライブを一つのチャネル上に接続してディスクドライブにアクセスする方法としては、ファイバチャネルアービトラレィッドループ(以下FC-A-Lと呼ぶ)がある。

[0028] 図10に、FC-A-Lの接続形態をディスクドライブDK1、DK2、DK3、DK4の接続形態を例として示す。各ディスクドライブの入出力ポートおよびスイッチSW1の入出力ポートは、送信機Txと受信機Rxを有する。FC-A-Lの接続形態は、例えば図10に示すように、各ドライブの入出力ポートおよびスイッチの入出力ポートをループ状に接続するポートロジである。各ドライブの入出力ポートはファイバチャネルのNL(Node Loop)ポートとして機能する。NLポートとは、ループ動作をする装置(ここではディスクアレイDA1接続側)出力ポートは、ファイバチャネルのFL(Fabric Loop)ポートとして機能する。FLポートとは、FC-A-Lを接続可能なスイッチのポートである。FLポートを有するループは、ファイバチャネルのパブリックループとして機能するので、

チャネルD11が形成するFC-Aはバプリックグループとなる。バプリックとは、ループ上のディスクドライブが、スイッチを介してループ外のポートと通信可能なループである。よって、ディスクドライブDK1、DK2、DK3、DK4は、スイッチSW1およびチャネルD01を介してディスクアダプタDKAと通信可能である。以上、チャネルD11の接続形態に照らしてみれば、チャネルD2、D13、D14でも同様である。本装置の形態のディスク装置においてRAIDシステムを構築する場合は、ディスクグループR1、R2、R3、R4を、それぞれRAIDグループとする。本装置の形態では、4個のディスクドライブでRAIDグループを構成しているが、RAIDグループを構成するドライブ数を4個に限るものではない。

【0032】 ディスクアレイDA1を構成するディスクドライブは、入出力ポートを5個有する。例えば、ディスクドライブDK1、DK2、DK3、DK4は、チャネルD1およびD21の高チャネルと接続する。ディスクアレイDA1は、チャネルD11とD21に接続した4個のディスクからなるディスクアレイと、D12とD22に接続した4個のディスクからなるディスクアレイと、D13とD23に接続した4個のディスクからなるディスクアレイと、D14とD24に接続した4個のディスクからなるディスクアレイ、からなる。チャネルD11、D12、D13、D14、D21、D22、D23、D24は、FC-ALでディスクドライブを接続する。

【0033】図11に本実施の形態におけるFC-A-
20 の接続形態を、ディスクドライブDK1、DK2、DK
3、DK4の接続形態を例として示す。各ディスクドラ
イブは、それぞれNLポートを2個有する。各ディスク
ドライブの入出力ポートおよびスイッチSW1、SW2
の入出力ポートは、送信端Txと受信端Rxを有する。

スイッチSW1、SW2のディスクレイドA1接続側
出力ポートは、FLポートである。チャネルD1に
より、スイッチSW1、ディスクドライブDK1、DK
2、DK3、DK4をループ状に接続する。同様にチャ
ネルD2により、スイッチSW2、ディスクドライブ
DK1、DK2、DK3、DK4をループ状に接続す
る。これら2個のループは、ファイバチャネルのパリ
ツクグループであり、ディスクドライブDK1、DK2、
DK3、DK4は、ファイバチャネルグループA1、

でディスクアダプタDKA1またはDKA2と通信可能である。以上、チャネルD11、D21の後続形態を例に説明したが、チャネルD12、D13、D14、D22、D23、D24でも同様である。本実施の形態のディスク装置においてRAIDシステムを構築する場合は、ディスクグループR1、R2、R3、R4を、それぞれRAIDグループとする。本実施の形態では、4個のディスクグループでRAIDグループを構成しているが、RAIDグループを構成するドライブ数を4個に限るものではない。

【0034】 ディスクアレイDA1内の全ディスクドライ
イブは、ディスクアダプタDKA1およびDKA2のど
ちからでもアクセス可能である。本実施の形態のディ
スク装置は、チャネルD1b、D2bをスイッチSW
1、SW2抜挿時の迂回経路として使用する。例えばス
イッチSW1が故障した場合でも、ディスクアダプタD
KA1はチャネルD1bとスイッチSW2経由でディスク
アレイDA1にアクセスできる。逆に、スイッチSW

2003 11 18 13:49

9 -

2が故障した場合は、ディスクアダプタDKA2はチャネルD2bとスイッチSW1経由でディスクアレイDA1にアクセスできるので、信頼性が高いディスク装置が実現できる。

【0035】図12に本発明の、第4の実施の形態で示す、本実施の形態のディスク装置は、第3の実施の形態のディスク装置に対して、スイッチSW1、SW2間を接続するチャネルD3 a、D3 bを設けた点に特徴がある。本実施の形態において、ディスクアダプタDKA1、DKA2と、スイッチSW1、SW2と、ディスクアレイDA1との間のデータ伝送方式は、ファイバチャネルを使用している。本実施の形態のディスク装置は、ディスクコントローラDKC1、DKC2と、スイッチSW1、SW2と、ディスクアレイDA1からなる。スイッチSW1とSW2は、第1の実施の形態と同様に選路切換と多重化、逆多重化を行うスイッチとして機能する。ディスクコントローラDKC1は、チャネルアダプタCHA1と、キャッシュメモリCM1と、ディスクアダプタDKA1からなる。ディスクコントローラDKC2は、チャネルアダプタCHA2と、キャッシュメモリCM2と、ディスクアダプタDKA2からなる。ディスクアダプタDKA1とスイッチSW1をチャネルD1 aで接続し、ディスクアダプタDKA2とスイッチSW2をチャネルD2 aで接続し、ディスクアダプタDKA1とスイッチSW2をチャネルD1 bで接続し、ディスクアダプタDKA2とスイッチSW1をチャネルD2 bで接続する。さらに、スイッチSW1とSW2をチャネルD3 a、D3 bで接続する。

【0036】ディスクアレイDA1を構成するディスクドライブは、出力ポートを2個有する。例えば、ディスクドライブDK1、DK2、DK3、DK4は、チャネルD1およびD2の両チャネルと接続する。ディスクアレイDA1は、チャネルD1とD21に接続した4個のディスクからなるディスクアレイと、D12とした4個のディスクからなるディスクアレイと、D12とD22に接続した4個のディスクからなるディスクアレイと、D13とD23に接続した4個のディスクからなるディスクアレイと、D14とD24に接続した4個のディスクからなるディスクアレイ、からなる。チャネルD11、D12、D13、D14、D21、D22、D23、D24は、図11に示するようにFC-ALでディスクドライブを接続する。ディスクアレイDA1内の全ディスクドライブは、ディスクアダプタDKA1およびDKA2のどちらからでもアクセス可能である。本実施の形態のディスク装置においてRAIDシステムを構築する場合は、ディスクグループR1、R2、R3、R4を、それぞれRAIDグループとする。本実施の形態では、4個のディスクドライブでRAIDグループを構成しているが、RAIDグループを構成するドライブ数を4個に限るものではない。

【0037】 ディスクアダプタDKA1、DKA2とデ

- 10 -

2003 11 18 13:49

8 となる。第3の実施の形態において、ディスクアダプ
タ-ディスクアレイ間スループットを4Gbpsにする
ためには、チャネルD1aおよびD2aのデータ伝送速
度を、それぞれ4Gbpsに高める必要がある。以上か
ら、本実施の形態によれば、ディスクアダプタ-スイッ
チ間のチャネル当りデータ伝送速度が低くても、ディ
スクアダプタ-ディスクアレイ間の総スループットが高い
ディスク装置が実現できる。

【0039】

【発明の効果】以上説明したように、本発明によれば以
下の効果がある。ディスクアダプタとディスクアレイ間
のスループットが高いディスク装置を提供できる。ま
た、ディスクアダプタとディスクアレイ間のスループッ
トが高く、且つディスクドライブ接続台数が多いディ
スク装置を提供できる。また、信頼性の高いディスクア
レイを有するディスク装置を提供できる。また、信頼性が
高いディスクアダプタとディスクアレイ間ネットワーク
を有するディスク装置を提供できる。また、信頼性およ
びスループットが高いディスクアダプタとディスクアレイ
間ネットワークを有するディスク装置を提供でき
る。また、ディスクからの読み出しおよびディスクへの
書き込みを高スループット化できるディスク装置を提供
できる。また、高スループットを維持できるディスク装
置を提供できる。また、ディスクアダプタとディスクア
レイ間のスループットが高く低コストなディスク装置を
提供できる。

【図面の簡単な説明】

【図1】本発明の第1の実施の形態のディスク装置を示
す図である。

【図2】本発明に用いるスイッチの構成を示す図であ
る。

【図3】本発明に用いるスイッチの構成を示す図であ
る。

【図4】本発明に用いるスイッチの動作を示す図であ
る。

【図5】本発明に用いるスイッチの動作を示す図であ
る。

【図6】本発明第1の実施の形態に対して、ディスクド
ライバを増設する方法を示す図である。

【図7】本発明の第2の実施の形態のディスク装置を示
す図である。

【図8】本発明の第3の実施の形態のディスク装置を示
す図である。

【図9】従来のディスク装置を示す図である。

【図10】FC-ALによる接続形態を説明する図であ
る。

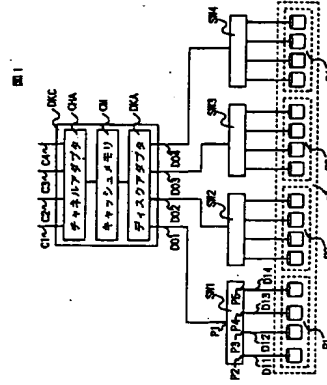
【図11】FC-ALによる接続形態を説明する図であ
る。

【図12】本発明の第4の実施の形態のディスク装置を
示す図である。

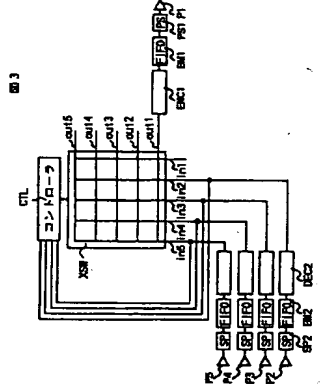
【符号の説明】

- DKC, DKC1, DKC2 ディスクコントローラ
- CHA, CHA1, CHA2 チャネルアダプタ
- CM, CM1, CM2 キャッシュメモリ
- DKA, DKA1, DKA2 ディスクアダプタ
- DA, DA1~DA4 ディスクアレイ
- DK1~DK4 ディスクドライブ
- R1~R6 ディスクグループ
- C1~C4, D1~D4, D01~D05, D11~D
14, D21~D24, D1a, D1b, D2a, D2
b, D3a, D3b チャネル
- SW1~SW4 スイッチ
- P1~P5 スイッチポート
- XSW クロスバススイッチ
- CTL スイッチコントローラ
- In1~In5 クロスバススイッチ入力ポート
- out1~out5 クロスバススイッチ出力ポート
- SP1, SP2 シリアルパラレル変換装置
- PS1, PS2 パラレルシリアル変換装置
- BM1, BM2 バッファメモリ
- DEC1, DEC2 8B10B変換デコーダ
- ENC1, ENC2 8B10B変換エンコーダ
- T1, T2, T3,
- T4, T5, T6 フレームの時間
- Tx 送信機
- Rx 受信機
- NL NLポート
- FL FLポート

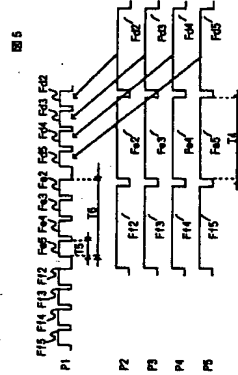
【図1】



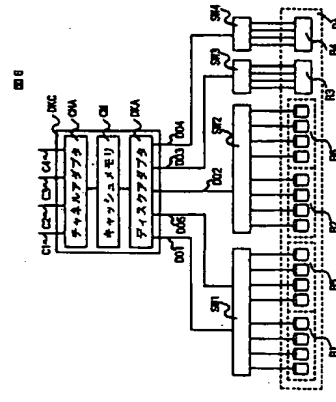
【図3】



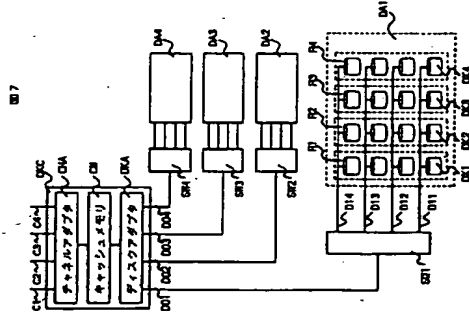
【図5】



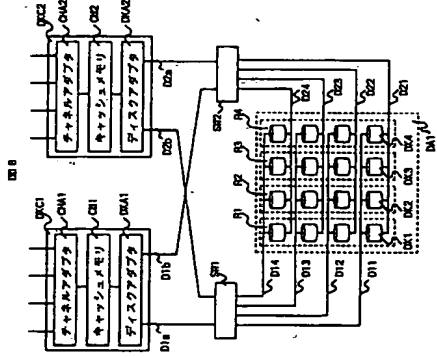
【図6】



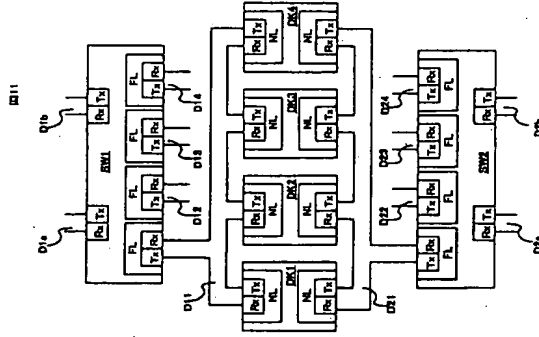
【図 7】



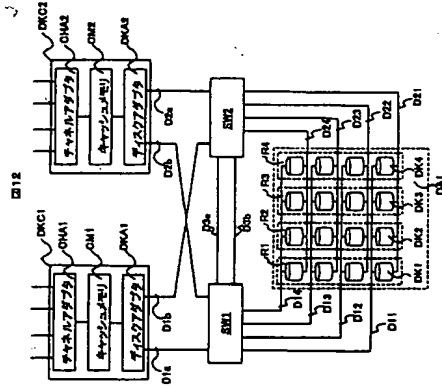
【図 8】



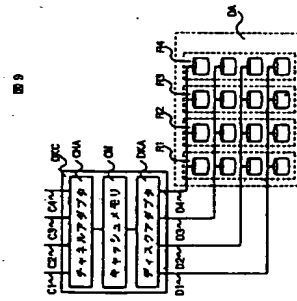
【図 11】



【図 12】



【図 9】



【図 10】

